

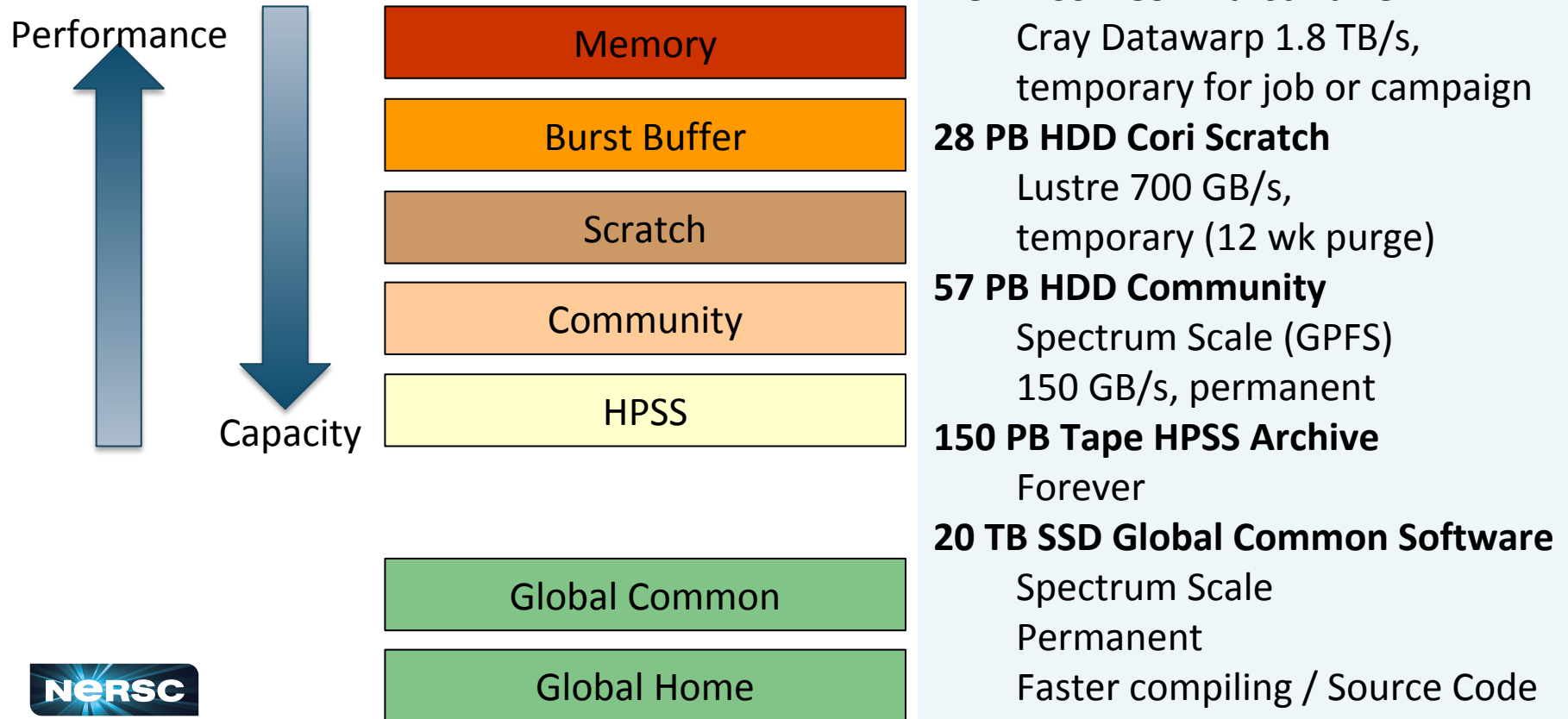
NERSC File Systems



New User Training
June 16, 2020

Wahid Bhimji
Data And Analytics SGroup

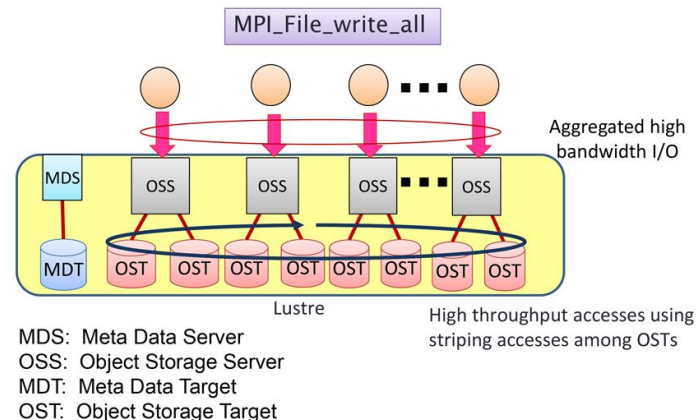
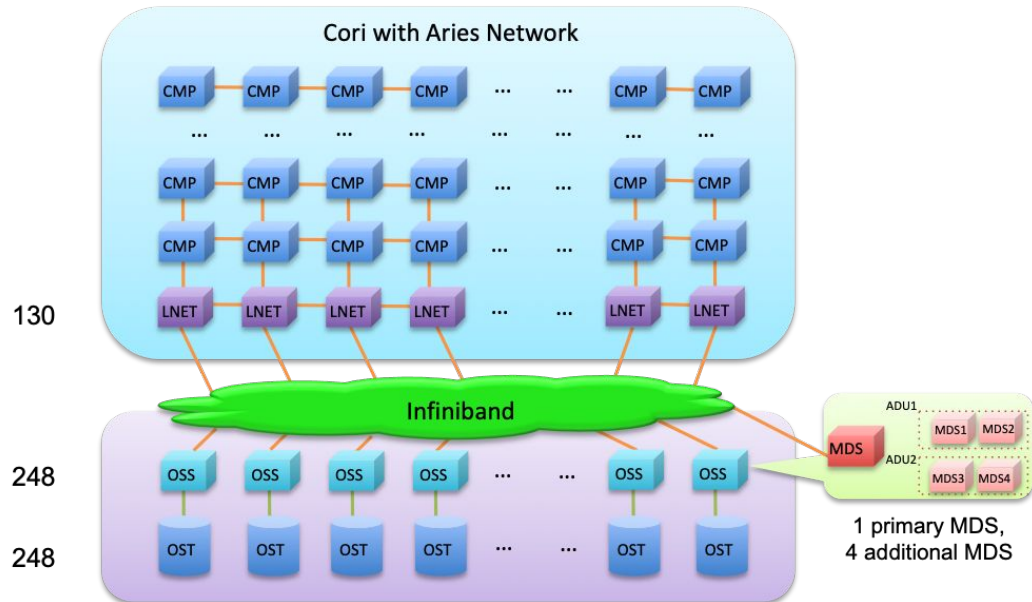
Simplified NERSC File Systems



Cori Scratch

Lustre, one of the most successful/mature HPC FS

Purged! Files not accessed for more than 12 weeks are automatically deleted



Using MPI-IO on Lustre[2]

“Scratch”: Optimize Performance with Striping

Scratch Striping Recommendations

- By default data on 1 OST, ideal for small files and file-per-process IO
- Single shared file IO should be striped according to its size
- Helper scripts

`stripe_small, stripe_medium stripe_large`

- Manually query with
`lfs getstripe <file_name>`

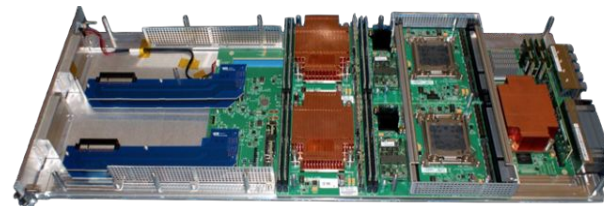
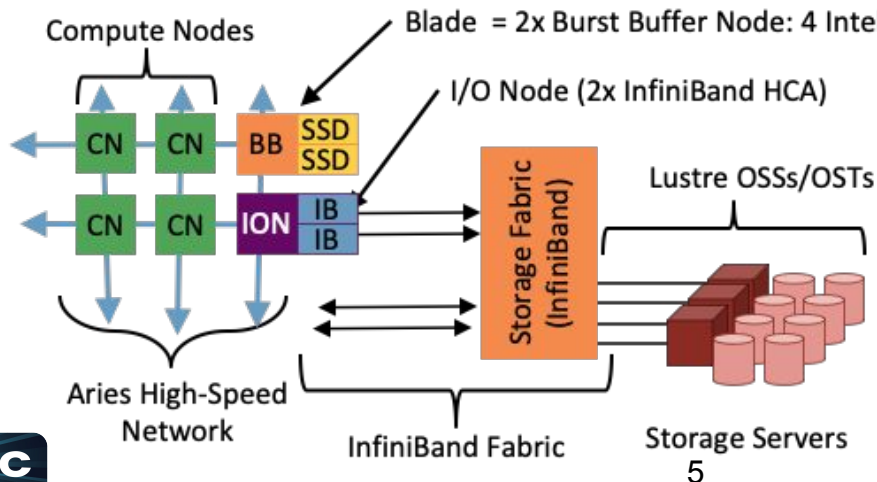
- Manually set with
`lfs setstripe -S 1m -c 2 <empty_folder>`

Single Shared-File I/O		File per Process
File size (GB)	command	
< 1	keep default striping	keep default striping
1 - 10	<code>stripe_small</code>	keep default striping
10 - 100	<code>stripe_medium</code>	keep default striping
> 100	<code>stripe_large</code>	keep default striping
> 1000	<code>stripe_large</code>	<code>stripe_large</code>

Burst Buffer (BB)

Datawarp (DW): Cray's applications I/O accelerator

- For: Data read in/out by high IO-bandwidth or IOPS application
- Transient: Allocated per-job or per-campaign ('persistent') via SLURM integration
- Users see mounted POSIX filesystem. Striped across BB nodes.



Burst Buffer Example

```
#!/bin/bash
#SBATCH -q regular -N 10 -C haswell -t 00:10:00
#DW jobdw capacity=1000GB access_mode=striped type=scratch
#DW stage_in source=/global/cscratch1/sd/username/file.dat destination=$DW_JOB_STRIPED/ type=file
#DW stage_out source=$DW_JOB_STRIPED/outputs destination=/lustre/outputs type=directory
srun my.x --infile=$DW_JOB_STRIPED/file.dat --outdir=$DW_JOB_STRIPED/outputs
```

- ‘jobdw’ – duration just for compute job (i.e. not ‘persistent’)
- ‘access_mode=striped’ – visible to all compute nodes, striped across BB nodes
 - Number of BB nodes depends on size requested - ‘granularity’ is 20 GB
- Data ‘stage_in’ before job start and ‘stage_out’ after
- \$DW_JOB_STRIPED env variable points to the mountpoint
- Can also use interactively:

```
wbhimji@cori12:~> cat bbf.conf
#DW jobdw capacity=1000GB access_mode=striped type=scratch
wbhimji@cori12:~> salloc -q interactive -N 1 -C knl --time=00:30:00 --bbf=bbf.conf
```


Creating a Persistent Reservation on Burst Buffer

Can be used by any job (set unix file permissions to share)

Don't forget to delete the PR when not needed (and no more than 6 weeks after)

DW not for long term storage and not resilient - stage_out anything you cannot afford to lose

Create a PR

```
#!/bin/bash
#SBATCH -q regular -N 10 -C haswell -t 00:10:00
#BB create_persistent name=myBBname
capacity=1000GB access=striped type=scratch
```

Delete PR

```
#!/bin/bash
#SBATCH -q regular -N 10 -C haswell
-t 00:10:00
#BB destroy_persistent name=myBBname
```

Can check reservation outside job:

```
wbhimji@cori12:~> scontrol show burst
Name=wahid_test_apr15_2
CreateTime=2020-02-14T16:10:36 Pool=(null)
Size=61872MiB State=allocated
UserID=wbhimji(68441)
```

Use PR in your jobs

```
#SBATCH -q regular -N 10 -C haswell -t 00:10:00
#DW peristentdw name=myBBname
#DW stage_in
source=/global/cscratch1/sd/[username]/inputs
destination=$DW_PERSISTENT_STRIPED_myBBname
type=directory
srun my.x
--indir=$DW_PERSISTENT_STRIPED_myBBname/inputs
```

Community File System

- For: Large datasets that you need for a longer period
- Set up for sharing with group read permissions by default
- Not for intensive I/O - use Scratch instead
- Can share data externally by dropping it into a www directory
Example: /global/cfs/cdirs/das/www/[username]
[https://portal.nersc.gov/project/das/\[username\]/](https://portal.nersc.gov/project/das/[username]/)
- Data is never purged. Snapshots. Usage is managed by quotas
- Projects can split their space allocations between multiple directories and give **separate** working groups **separate** quotas
 - Environment variable \$CFS points to /global/cfs/cdirs

<https://docs.nersc.gov/filesystems/community/>

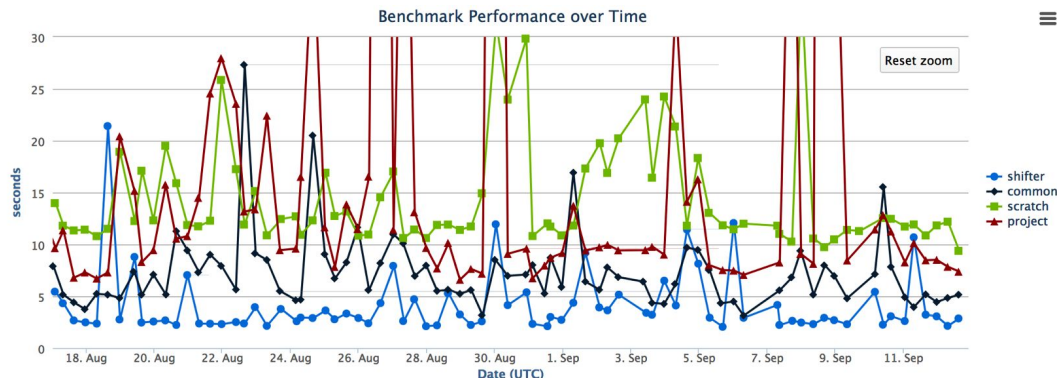
HPSS

- For: Data from your finished paper, raw data you might need in case of emergency, really hard to generate data
- HPSS is tape!
 - Data first hits a spinning disk cache and gets migrated to tapes
 - Files can end up spread all over, so use htar to aggregate into bundles of 100 GB - 2 TB
 - Archive the way you intend to retrieve the data

<https://docs.nersc.gov/filesystems/archive/>
https://docs.nersc.gov/filesystems/archive_access/

“Global Common”: Software Filesystem

- For: Software stacks - Why? Library load performance

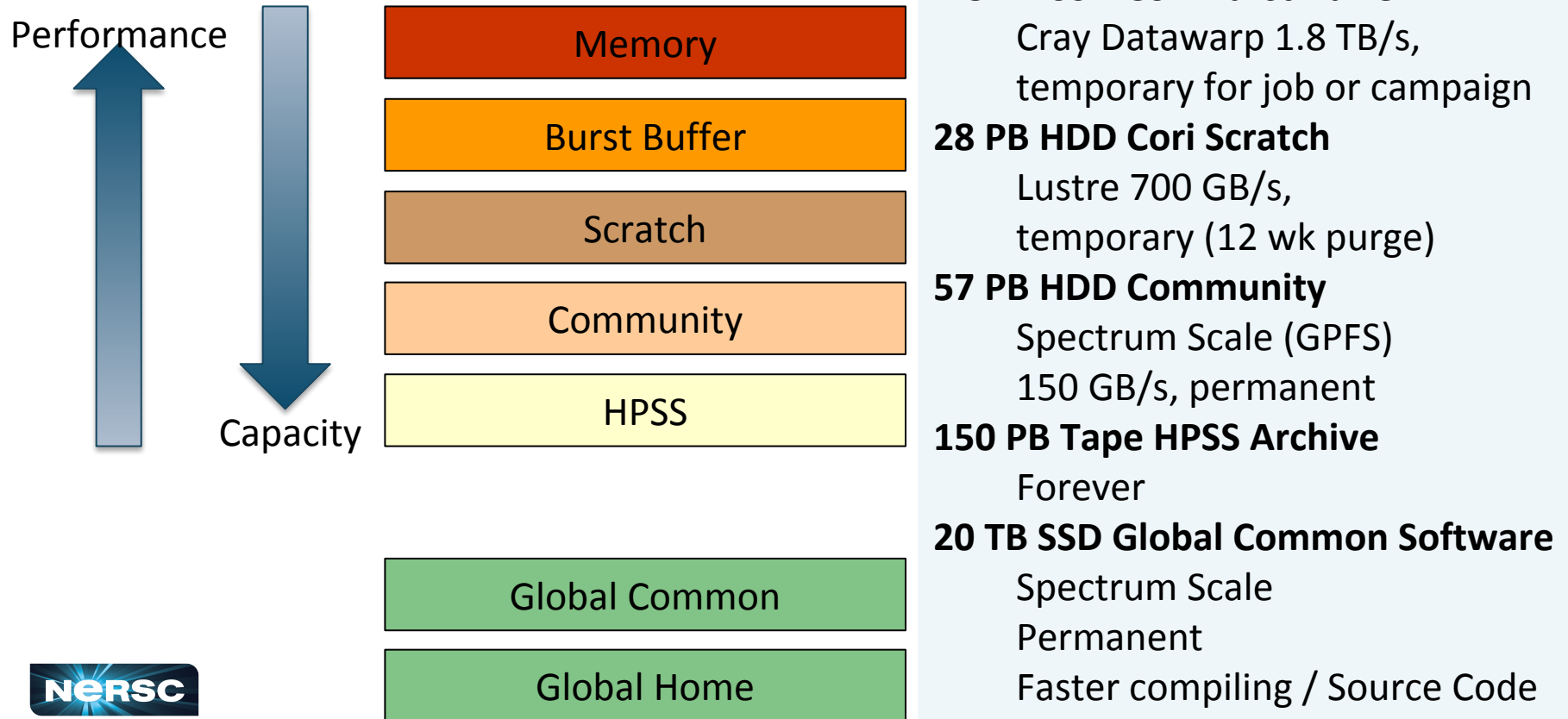


- Group writable directories similar to community, but with a smaller quota, /global/common/software/<projectname>
 - Write from login node; read-only on compute node
- Smaller block size for faster compiles than project

Home Directories

- For: Source files, compiling, configuration files
- 20G quota
- Not intended for intensive I/O (e.g. application I/O) - use Scratch instead
- Backed up monthly by HPSS
- Snapshots are also available e.g. my homedir is at `/global/homes/.snapshots/2020-06-14/w/wbhimji`

Simplified NERSC File Systems



Data Dashboard

Data Dashboard

Showing disk space and inode usage for project directories at NERSC to which you have access as PI, PI proxy, or user (includes /project, /p

CAL directory in /project

[Toggle Usage Details](#)[My Files and Dirs](#)[Browse](#)

carver directory in /project

[Toggle Usage Details](#)[My Files and Dirs](#)[Browse](#)

coe directory in /project

[Toggle Usage Details](#)[My Files and Dirs](#)[Browse](#)

cosmo directory in /project

[Toggle Usage Details](#)[My Files and Dirs](#)[Browse](#)

crcns directory in /project

[Toggle Usage Details](#)[My Files and Dirs](#)[Browse](#)

das directory in /project

[Toggle Usage Details](#)[My Files and Dirs](#)[Browse](#)

Data Dashboard: Usage Reports

Data Dashboard

Showing disk space and inode usage for project directories at NERSC to which you have access as PI, PI proxy, or user (includes /project, /projecta, and /projectb/sandbox)

CAL directory in /project

[Toggle Usage Details](#)

[My Files and Dirs](#)

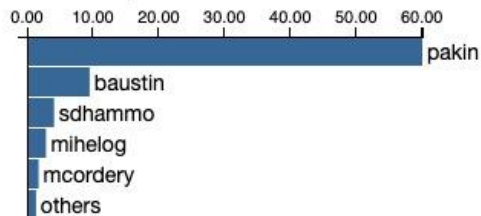
[Browse](#)



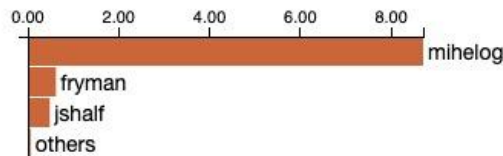
Data as of Thu Jun 20 2019 23:59:59 GMT-0700 (Pacific Daylight Time)

Breakdown of allocation usage:

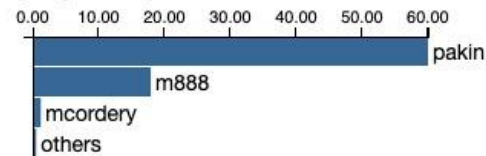
user % of space allocation



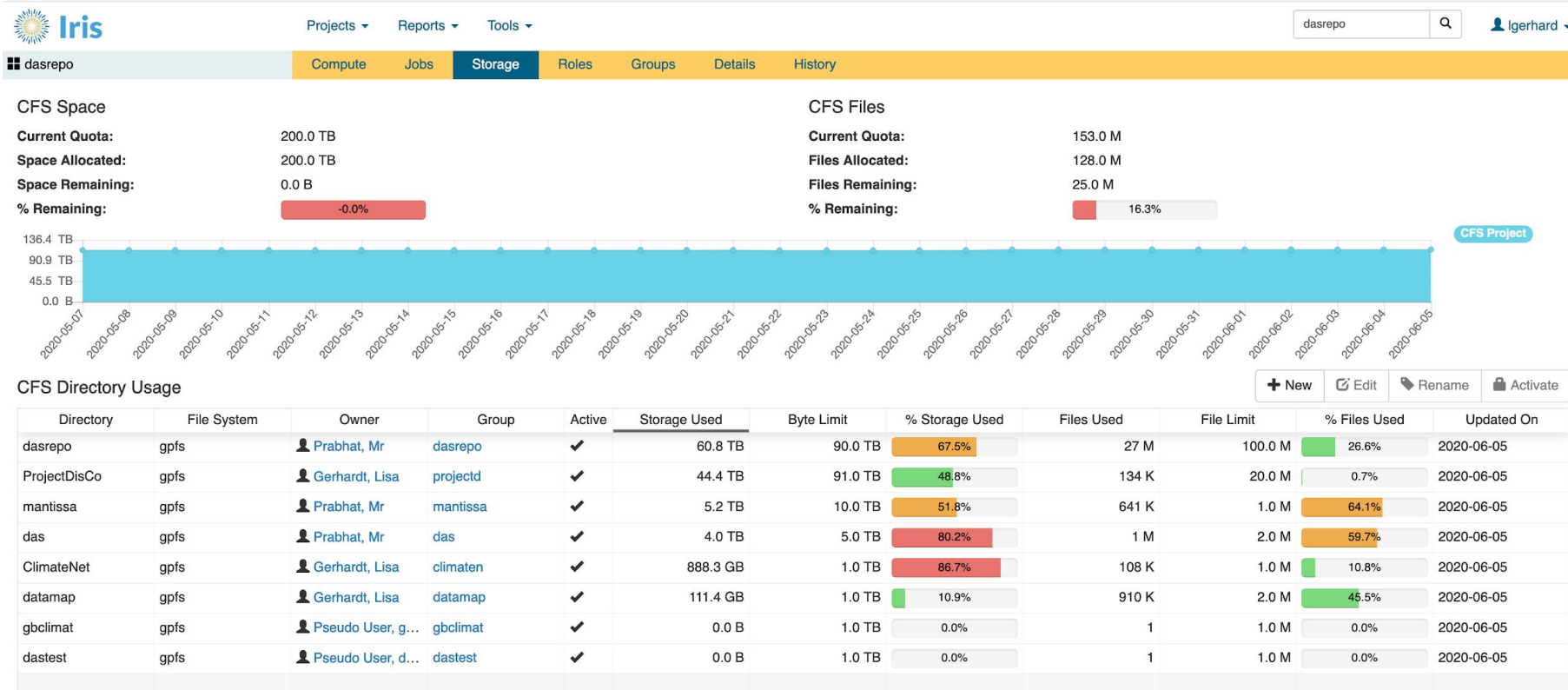
user % of inode allocation



group % of space allocation



Adjusting Quotas in IRIS



Resources

- Cori File Systems

<https://docs.nersc.gov/filesystems//>

- NERSC Burst Buffer Web Pages

<https://docs.nersc.gov/filesystems/cori-burst-buffer/>

- Example batch scripts

<https://docs.nersc.gov/jobs/examples/#burst-buffer>

- DataWarp Users Guide

https://pubs.cray.com/bundle/XC_Series_DataWarp_User_Guide_S-2558_public_S2558_final/page/About_XC_Series_DataWarp_User_Guide.html

Thank You and
Welcome to
NERSC!

